```
In [1]: import numpy as np
        import pandas as pd
```

```
In [2]: movies = pd.read_csv('tmdb_5000_movies.csv')
        credits = pd.read_csv('tmdb_5000_credits.csv')
```

## Data pre-processing

```
In [3]: movies.head()
```

Out[3]:

| | budget | genres | homepage | id | keywords | original_language | original_title | overview | popularity | production_companies | production_countries | release_date | revenue | runtime | spoken_languages | status | tagline | title | vote_average | vote_count |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 237000000 | [{"id": 28, "name": "Action"}, {"id": 12, "nam... | http://www.avatarmovie.com/ | 19995 | [{"id": 1463, "name": "culture clash"}, {"id":... | en | Avatar | In the 22nd century, a paraplegic Marine is di... | 150.437577 | [{"name": "Ingenious Film Partners", "id": 289... | [{"iso_3166_1": "US", "name": "United States o... | 2009-12-10 | 2787965087 | 162.0 | [{"iso_639_1": "en", "name": "English"}, {"iso... | Released | Enter the World of Pandora. | Avatar | 7.2 | 11800 |
| 1 | 300000000 | [{"id": 12, "name": "Adventure"}, {"id": 14, "... | http://disney.go.com/disneypictures/pirates/ | 285 | [{"id": 270, "name": "ocean"}, {"id": 726, "na... | en | Pirates of the Caribbean: At World's End | Captain Barbossa, long believed to be dead, ha... | 139.082615 | [{"name": "Walt Disney Pictures", "id": 2}, {"... | [{"iso_3166_1": "US", "name": "United States o... | 2007-05-19 | 961000000 | 169.0 | [{"iso_639_1": "en", "name": "English"}] | Released | At the end of the world, the adventure begins. | Pirates of the Caribbean: At World's End | 6.9 | 4500 |
| 2 | 245000000 | [{"id": 28, "name": "Action"}, {"id": 12, "nam... | http://www.sonypictures.com/movies/spectre/ | 206647 | [{"id": 470, "name": "spy"}, {"id": 818, "name... | en | Spectre | A cryptic message from Bond's past sends him o... | 107.376788 | [{"name": "Columbia Pictures", "id": 5}, {"nam... | [{"iso_3166_1": "GB", "name": "United Kingdom"... | 2015-10-26 | 880674609 | 148.0 | [{"iso_639_1": "fr", "name": "Fran\u00e7ais"},... | Released | A Plan No One Escapes | Spectre | 6.3 | 4466 |
| 3 | 250000000 | [{"id": 28, "name": "Action"}, {"id": 80, "nam... | http://www.thedarkknightrises.com/ | 49026 | [{"id": 849, "name": "dc comics"}, {"id": 853,... | en | The Dark Knight Rises | Following the death of District Attorney Harve... | 112.312950 | [{"name": "Legendary Pictures", "id": 923}, {"... | [{"iso_3166_1": "US", "name": "United States o... | 2012-07-16 | 1084939099 | 165.0 | [{"iso_639_1": "en", "name": "English"}] | Released | The Legend Ends | The Dark Knight Rises | 7.6 | 9106 |
| 4 | 260000000 | [{"id": 28, "name": "Action"}, {"id": 12, "nam... | http://movies.disney.com/john-carter | 49529 | [{"id": 818, "name": "based on novel"}, {"id":... | en | John Carter | John Carter is a war-weary, former military ca... | 43.926995 | [{"name": "Walt Disney Pictures", "id": 2}] | [{"iso_3166_1": "US", "name": "United States o... | 2012-03-07 | 284139100 | 132.0 | [{"iso_639_1": "en", "name": "English"}] | Released | Lost in our world, found in another. | John Carter | 6.1 | 2124 |

```
In [4]: credits.head()
```

Out[4]:

| | movie_id | title | cast | crew |
|---|---|---|---|---|
| 0 | 19995 | Avatar | [{"cast_id": 242, "character": "Jake Sully", "... | [{"credit_id": "52fe48009251416c750aca23", "de... |
| 1 | 285 | Pirates of the Caribbean: At World's End | [{"cast_id": 4, "character": "Captain Jack Spa... | [{"credit_id": "52fe4232c3a36847f800b579", "de... |
| 2 | 206647 | Spectre | [{"cast_id": 1, "character": "James Bond", "cr... | [{"credit_id": "54805967c3a36829b5002c41", "de... |
| 3 | 49026 | The Dark Knight Rises | [{"cast_id": 2, "character": "Bruce Wayne / Ba... | [{"credit_id": "52fe4781c3a36847f81398c3", "de... |
| 4 | 49529 | John Carter | [{"cast_id": 5, "character": "John Carter", "c... | [{"credit_id": "52fe479ac3a36847f813eaa3", "de... |

```
In [5]: movies = movies.merge(credits,on="title")
```

```
In [6]: movies.shape
```

Out[6]: (4809, 23)

```
In [7]: movies.head(1)
```

Out[7]:

| | budget | genres | homepage | id | keywords | original_language | original_title | overview | popularity | production_companies | ... | runtime | spoken_languages | status | tagline | title | vote_average | vote_count | movie_id | cast | crew |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 237000000 | [{"id": 28, "name": "Action"}, {"id": 12, "nam... | http://www.avatarmovie.com/ | 19995 | [{"id": 1463, "name": "culture clash"}, {"id":... | en | Avatar | In the 22nd century, a paraplegic Marine is di... | 150.437577 | [{"name": "Ingenious Film Partners", "id": 289... | ... | 162.0 | [{"iso_639_1": "en", "name": "English"}, {"iso... | Released | Enter the World of Pandora. | Avatar | 7.2 | 11800 | 19995 | [{"cast_id": 242, "character": "Jake Sully", "... | [{"credit_id": "52fe48009251416c750aca23", "de... |

1 rows × 23 columns

```
In [8]: # genres #id #keywords #title # overview #cast #crew

        movies = movies[['movie_id','title','overview','genres','keywords','cast','crew']]
```

In [9]: `movies.info()`

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 4809 entries, 0 to 4808
Data columns (total 7 columns):
 #   Column    Non-Null Count  Dtype
---  ------    --------------  -----
 0   movie_id  4809 non-null   int64
 1   title     4809 non-null   object
 2   overview  4806 non-null   object
 3   genres    4809 non-null   object
 4   keywords  4809 non-null   object
 5   cast      4809 non-null   object
 6   crew      4809 non-null   object
dtypes: int64(1), object(6)
memory usage: 300.6+ KB
```

In [10]: `movies`

Out[10]:

| | movie_id | title | overview | genres | keywords | cast | crew |
|---|---|---|---|---|---|---|---|
| 0 | 19995 | Avatar | In the 22nd century, a paraplegic Marine is di... | [{"id": 28, "name": "Action"}, {"id": 12, "nam... | [{"id": 1463, "name": "culture clash"}, {"id":... | [{"cast_id": 242, "character": "Jake Sully", "... | [{"credit_id": "52fe48009251416c750aca23", "de... |
| 1 | 285 | Pirates of the Caribbean: At World's End | Captain Barbossa, long believed to be dead, ha... | [{"id": 12, "name": "Adventure"}, {"id": 14, "... | [{"id": 270, "name": "ocean"}, {"id": 726, "na... | [{"cast_id": 4, "character": "Captain Jack Spa... | [{"credit_id": "52fe4232c3a36847f800b579", "de... |
| 2 | 206647 | Spectre | A cryptic message from Bond's past sends him o... | [{"id": 28, "name": "Action"}, {"id": 12, "nam... | [{"id": 470, "name": "spy"}, {"id": 818, "name... | [{"cast_id": 1, "character": "James Bond", "cr... | [{"credit_id": "54805967c3a36829b5002c41", "de... |
| 3 | 49026 | The Dark Knight Rises | Following the death of District Attorney Harve... | [{"id": 28, "name": "Action"}, {"id": 80, "nam... | [{"id": 849, "name": "dc comics"}, {"id": 853,... | [{"cast_id": 2, "character": "Bruce Wayne / Ba... | [{"credit_id": "52fe4781c3a36847f81398c3", "de... |
| 4 | 49529 | John Carter | John Carter is a war-weary, former military ca... | [{"id": 28, "name": "Action"}, {"id": 12, "nam... | [{"id": 818, "name": "based on novel"}, {"id":... | [{"cast_id": 5, "character": "John Carter", "c... | [{"credit_id": "52fe479ac3a36847f813eaa3", "de... |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 4804 | 9367 | El Mariachi | El Mariachi just wants to play his guitar and ... | [{"id": 28, "name": "Action"}, {"id": 80, "nam... | [{"id": 5616, "name": "united states\u2013mexi... | [{"cast_id": 1, "character": "El Mariachi", "c... | [{"credit_id": "52fe44eec3a36847f80b280b", "de... |
| 4805 | 72766 | Newlyweds | A newlywed couple's honeymoon is upended by th... | [{"id": 35, "name": "Comedy"}, {"id": 10749, "... | [] | [{"cast_id": 1, "character": "Buzzy", "credit_... | [{"credit_id": "52fe487dc3a368484e0fb013", "de... |
| 4806 | 231617 | Signed, Sealed, Delivered | "Signed, Sealed, Delivered" introduces a dedic... | [{"id": 35, "name": "Comedy"}, {"id": 18, "nam... | [{"id": 248, "name": "date"}, {"id": 699, "nam... | [{"cast_id": 8, "character": "Oliver O\u2019To... | [{"credit_id": "52fe4df3c3a36847f8275ecf", "de... |
| 4807 | 126186 | Shanghai Calling | When ambitious New York attorney Sam is sent t... | [] | [] | [{"cast_id": 3, "character": "Sam", "credit_id... | [{"credit_id": "52fe4ad9c3a368484e16a36b", "de... |
| 4808 | 25975 | My Date with Drew | Ever since the second grade when he first saw ... | [{"id": 99, "name": "Documentary"}] | [{"id": 1523, "name": "obsession"}, {"id": 224... | [{"cast_id": 3, "character": "Herself", "credi... | [{"credit_id": "58ce021b9251415a390165d9", "de... |

4809 rows × 7 columns

In [11]: `movies.isnull().sum()`

Out[11]:
```
movie_id    0
title       0
overview    3
genres      0
keywords    0
cast        0
crew        0
dtype: int64
```

In [12]: `movies.dropna(inplace=True)`

In [13]: `movies.iloc[0].genres`

Out[13]: `'[{"id": 28, "name": "Action"}, {"id": 12, "name": "Adventure"}, {"id": 14, "name": "Fantasy"}, {"id": 878, "name": "Science Fiction"}]'`

In [14]:
```python
def convert(obj):
    L = []
    for i in ast.literal_eval(obj):
        L.append(i['name'])
    return L
```

In [15]: `import ast`

In [16]: `movies['genres'] = movies['genres'].apply(convert)`

In [17]: `movies.head(2)`

Out[17]:

| | movie_id | title | overview | genres | keywords | cast | crew |
|---|---|---|---|---|---|---|---|
| 0 | 19995 | Avatar | In the 22nd century, a paraplegic Marine is di... | [Action, Adventure, Fantasy, Science Fiction] | [{"id": 1463, "name": "culture clash"}, {"id":... | [{"cast_id": 242, "character": "Jake Sully", "... | [{"credit_id": "52fe48009251416c750aca23", "de... |
| 1 | 285 | Pirates of the Caribbean: At World's End | Captain Barbossa, long believed to be dead, ha... | [Adventure, Fantasy, Action] | [{"id": 270, "name": "ocean"}, {"id": 726, "na... | [{"cast_id": 4, "character": "Captain Jack Spa... | [{"credit_id": "52fe4232c3a36847f800b579", "de... |

In [18]: `movies['keywords'] = movies['keywords'].apply(convert)`

In [19]: `movies.head(2)`

Out[19]:

| | movie_id | title | overview | genres | keywords | cast | crew |
|---|---|---|---|---|---|---|---|
| 0 | 19995 | Avatar | In the 22nd century, a paraplegic Marine is di... | [Action, Adventure, Fantasy, Science Fiction] | [culture clash, future, space war, space colon... | [{"cast_id": 242, "character": "Jake Sully", "... | [{"credit_id": "52fe48009251416c750aca23", "de... |
| 1 | 285 | Pirates of the Caribbean: At World's End | Captain Barbossa, long believed to be dead, ha... | [Adventure, Fantasy, Action] | [ocean, drug abuse, exotic island, east india ... | [{"cast_id": 4, "character": "Captain Jack Spa... | [{"credit_id": "52fe4232c3a36847f800b579", "de... |

In [20]:
```python
def convert3(obj):
    L = []
    counter = 0
    for i in ast.literal_eval(obj):
        if counter != 3:
            L.append(i['name'])
            counter+=1
        else:
            break
    return L
```

In [21]: `movies['cast'] = movies['cast'].apply(convert3)`

In [22]: `movies.head()`

Out[22]:

| | movie_id | title | overview | genres | keywords | cast | crew |
|---|---|---|---|---|---|---|---|
| 0 | 19995 | Avatar | In the 22nd century, a paraplegic Marine is di... | [Action, Adventure, Fantasy, Science Fiction] | [culture clash, future, space war, space colon... | [Sam Worthington, Zoe Saldana, Sigourney Weaver] | [{"credit_id": "52fe48009251416c750aca23", "de... |
| 1 | 285 | Pirates of the Caribbean: At World's End | Captain Barbossa, long believed to be dead, ha... | [Adventure, Fantasy, Action] | [ocean, drug abuse, exotic island, east india ... | [Johnny Depp, Orlando Bloom, Keira Knightley] | [{"credit_id": "52fe4232c3a36847f800b579", "de... |
| 2 | 206647 | Spectre | A cryptic message from Bond's past sends him o... | [Action, Adventure, Crime] | [spy, based on novel, secret agent, sequel, mi... | [Daniel Craig, Christoph Waltz, Léa Seydoux] | [{"credit_id": "54805967c3a36829b5002c41", "de... |
| 3 | 49026 | The Dark Knight Rises | Following the death of District Attorney Harve... | [Action, Crime, Drama, Thriller] | [dc comics, crime fighter, terrorist, secret i... | [Christian Bale, Michael Caine, Gary Oldman] | [{"credit_id": "52fe4781c3a36847f81398c3", "de... |
| 4 | 49529 | John Carter | John Carter is a war-weary, former military ca... | [Action, Adventure, Science Fiction] | [based on novel, mars, medallion, space travel... | [Taylor Kitsch, Lynn Collins, Samantha Morton] | [{"credit_id": "52fe479ac3a36847f813eaa3", "de... |

In [23]:
```python
def fetch_director(obj):
    L = []
    for i in ast.literal_eval(obj):
        if i['job'] == 'Director':
            L.append(i['name'])
            break
    return L
```

In [24]: `movies['crew'] = movies['crew'].apply(fetch_director)`

In [25]: `movies.head()`

Out[25]:

| | movie_id | title | overview | genres | keywords | cast | crew |
|---|---|---|---|---|---|---|---|
| 0 | 19995 | Avatar | In the 22nd century, a paraplegic Marine is di... | [Action, Adventure, Fantasy, Science Fiction] | [culture clash, future, space war, space colon... | [Sam Worthington, Zoe Saldana, Sigourney Weaver] | [James Cameron] |
| 1 | 285 | Pirates of the Caribbean: At World's End | Captain Barbossa, long believed to be dead, ha... | [Adventure, Fantasy, Action] | [ocean, drug abuse, exotic island, east india ... | [Johnny Depp, Orlando Bloom, Keira Knightley] | [Gore Verbinski] |
| 2 | 206647 | Spectre | A cryptic message from Bond's past sends him o... | [Action, Adventure, Crime] | [spy, based on novel, secret agent, sequel, mi... | [Daniel Craig, Christoph Waltz, Léa Seydoux] | [Sam Mendes] |
| 3 | 49026 | The Dark Knight Rises | Following the death of District Attorney Harve... | [Action, Crime, Drama, Thriller] | [dc comics, crime fighter, terrorist, secret i... | [Christian Bale, Michael Caine, Gary Oldman] | [Christopher Nolan] |
| 4 | 49529 | John Carter | John Carter is a war-weary, former military ca... | [Action, Adventure, Science Fiction] | [based on novel, mars, medallion, space travel... | [Taylor Kitsch, Lynn Collins, Samantha Morton] | [Andrew Stanton] |

In [26]: `movies['overview'][0]`

Out[26]: `'In the 22nd century, a paraplegic Marine is dispatched to the moon Pandora on a unique mission, but becomes torn between following orders and protecting an alien civilization.'`

In [27]: `movies['overview'] = movies['overview'].apply(lambda x:x.split())`

In [28]: `movies.head()`

Out[28]:

| | movie_id | title | overview | genres | keywords | cast | crew |
|---|---|---|---|---|---|---|---|
| 0 | 19995 | Avatar | [In, the, 22nd, century,, a, paraplegic, Marin... | [Action, Adventure, Fantasy, Science Fiction] | [culture clash, future, space war, space colon... | [Sam Worthington, Zoe Saldana, Sigourney Weaver] | [James Cameron] |
| 1 | 285 | Pirates of the Caribbean: At World's End | [Captain, Barbossa,, long, believed, to, be, d... | [Adventure, Fantasy, Action] | [ocean, drug abuse, exotic island, east india ... | [Johnny Depp, Orlando Bloom, Keira Knightley] | [Gore Verbinski] |
| 2 | 206647 | Spectre | [A, cryptic, message, from, Bond's, past, send... | [Action, Adventure, Crime] | [spy, based on novel, secret agent, sequel, mi... | [Daniel Craig, Christoph Waltz, Léa Seydoux] | [Sam Mendes] |
| 3 | 49026 | The Dark Knight Rises | [Following, the, death, of, District, Attorney... | [Action, Crime, Drama, Thriller] | [dc comics, crime fighter, terrorist, secret i... | [Christian Bale, Michael Caine, Gary Oldman] | [Christopher Nolan] |
| 4 | 49529 | John Carter | [John, Carter, is, a, war-weary,, former, mili... | [Action, Adventure, Science Fiction] | [based on novel, mars, medallion, space travel... | [Taylor Kitsch, Lynn Collins, Samantha Morton] | [Andrew Stanton] |

```
In [29]:  movies['genres'] = movies['genres'].apply(lambda x:[i.replace(" ","")for i in x])
          movies['keywords'] = movies['keywords'].apply(lambda x:[i.replace(" ","")for i in x])
          movies['cast'] = movies['cast'].apply(lambda x:[i.replace(" ","")for i in x])
          movies['crew'] = movies['crew'].apply(lambda x:[i.replace(" ","")for i in x])
```

```
In [30]:  movies.head()
```

Out[30]:

| | movie_id | title | overview | genres | keywords | cast | crew |
|---|---|---|---|---|---|---|---|
| **0** | 19995 | Avatar | [In, the, 22nd, century,, a, paraplegic, Marin... | [Action, Adventure, Fantasy, ScienceFiction] | [cultureclash, future, spacewar, spacecolony, ... | [SamWorthington, ZoeSaldana, SigourneyWeaver] | [JamesCameron] |
| **1** | 285 | Pirates of the Caribbean: At World's End | [Captain, Barbossa,, long, believed, to, be, d... | [Adventure, Fantasy, Action] | [ocean, drugabuse, exoticisland, eastindiatrad... | [JohnnyDepp, OrlandoBloom, KeiraKnightley] | [GoreVerbinski] |
| **2** | 206647 | Spectre | [A, cryptic, message, from, Bond's, past, send... | [Action, Adventure, Crime] | [spy, basedonnovel, secretagent, sequel, mi6, ... | [DanielCraig, ChristophWaltz, LéaSeydoux] | [SamMendes] |
| **3** | 49026 | The Dark Knight Rises | [Following, the, death, of, District, Attorney... | [Action, Crime, Drama, Thriller] | [dccomics, crimefighter, terrorist, secretiden... | [ChristianBale, MichaelCaine, GaryOldman] | [ChristopherNolan] |
| **4** | 49529 | John Carter | [John, Carter, is, a, war-weary,, former, mili... | [Action, Adventure, ScienceFiction] | [basedonnovel, mars, medallion, spacetravel, p... | [TaylorKitsch, LynnCollins, SamanthaMorton] | [AndrewStanton] |

```
In [31]:  movies['tags'] = movies['overview'] + movies['genres'] + movies['keywords'] + movies['cast'] + movies['crew']
```

```
In [32]:  movies.head()
```

Out[32]:

| | movie_id | title | overview | genres | keywords | cast | crew | tags |
|---|---|---|---|---|---|---|---|---|
| **0** | 19995 | Avatar | [In, the, 22nd, century,, a, paraplegic, Marin... | [Action, Adventure, Fantasy, ScienceFiction] | [cultureclash, future, spacewar, spacecolony, ... | [SamWorthington, ZoeSaldana, SigourneyWeaver] | [JamesCameron] | [In, the, 22nd, century,, a, paraplegic, Marin... |
| **1** | 285 | Pirates of the Caribbean: At World's End | [Captain, Barbossa,, long, believed, to, be, d... | [Adventure, Fantasy, Action] | [ocean, drugabuse, exoticisland, eastindiatrad... | [JohnnyDepp, OrlandoBloom, KeiraKnightley] | [GoreVerbinski] | [Captain, Barbossa,, long, believed, to, be, d... |
| **2** | 206647 | Spectre | [A, cryptic, message, from, Bond's, past, send... | [Action, Adventure, Crime] | [spy, basedonnovel, secretagent, sequel, mi6, ... | [DanielCraig, ChristophWaltz, LéaSeydoux] | [SamMendes] | [A, cryptic, message, from, Bond's, past, send... |
| **3** | 49026 | The Dark Knight Rises | [Following, the, death, of, District, Attorney... | [Action, Crime, Drama, Thriller] | [dccomics, crimefighter, terrorist, secretiden... | [ChristianBale, MichaelCaine, GaryOldman] | [ChristopherNolan] | [Following, the, death, of, District, Attorney... |
| **4** | 49529 | John Carter | [John, Carter, is, a, war-weary,, former, mili... | [Action, Adventure, ScienceFiction] | [basedonnovel, mars, medallion, spacetravel, p... | [TaylorKitsch, LynnCollins, SamanthaMorton] | [AndrewStanton] | [John, Carter, is, a, war-weary,, former, mili... |

```
In [33]:  new_df = movies[['movie_id','title','tags']]
```

```
In [34]:  new_df
```

Out[34]:

| | movie_id | title | tags |
|---|---|---|---|
| **0** | 19995 | Avatar | [In, the, 22nd, century,, a, paraplegic, Marin... |
| **1** | 285 | Pirates of the Caribbean: At World's End | [Captain, Barbossa,, long, believed, to, be, d... |
| **2** | 206647 | Spectre | [A, cryptic, message, from, Bond's, past, send... |
| **3** | 49026 | The Dark Knight Rises | [Following, the, death, of, District, Attorney... |
| **4** | 49529 | John Carter | [John, Carter, is, a, war-weary,, former, mili... |
| **...** | ... | ... | ... |
| **4804** | 9367 | El Mariachi | [El, Mariachi, just, wants, to, play, his, gui... |
| **4805** | 72766 | Newlyweds | [A, newlywed, couple's, honeymoon, is, upended... |
| **4806** | 231617 | Signed, Sealed, Delivered | ["Signed,, Sealed,, Delivered", introduces, a,... |
| **4807** | 126186 | Shanghai Calling | [When, ambitious, New, York, attorney, Sam, is... |
| **4808** | 25975 | My Date with Drew | [Ever, since, the, second, grade, when, he, fi... |

4806 rows × 3 columns

```
In [35]:  new_df["tags"] = new_df["tags"].apply(lambda x:" ".join(x))
```

```
<ipython-input-35-cffb23d2b53a>:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)
  new_df["tags"] = new_df["tags"].apply(lambda x:" ".join(x))
```

```
In [36]:  new_df.head()
```

Out[36]:

| | movie_id | title | tags |
|---|---|---|---|
| **0** | 19995 | Avatar | In the 22nd century a paraplegic Marine is di... |
| **1** | 285 | Pirates of the Caribbean: At World's End | Captain Barbossa, long believed to be dead, ha... |
| **2** | 206647 | Spectre | A cryptic message from Bond's past sends him o... |
| **3** | 49026 | The Dark Knight Rises | Following the death of District Attorney Harve... |
| **4** | 49529 | John Carter | John Carter is a war-weary, former military ca... |

In [37]:
```python
new_df["tags"] = new_df["tags"].apply(lambda x:x.lower())
```

```
<ipython-input-37-3af545d32fd5>:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)
  new_df["tags"] = new_df["tags"].apply(lambda x:x.lower())
```

In [38]:
```python
new_df.head()
```

Out[38]:

| | movie_id | title | tags |
|---|---|---|---|
| 0 | 19995 | Avatar | in the 22nd century, a paraplegic marine is di... |
| 1 | 285 | Pirates of the Caribbean: At World's End | captain barbossa, long believed to be dead, ha... |
| 2 | 206647 | Spectre | a cryptic message from bond's past sends him o... |
| 3 | 49026 | The Dark Knight Rises | following the death of district attorney harve... |
| 4 | 49529 | John Carter | john carter is a war-weary, former military ca... |

In [39]:
```python
import nltk
```

In [40]:
```python
##pip install nltk
```

In [41]:
```python
from nltk.stem import PorterStemmer

ps = PorterStemmer()

def stem(text):
    y = []

    for i in text.split():
        y.append(ps.stem(i))

    return " ".join(y)

new_df['tags'] = new_df['tags'].apply(stem)
```

```
<ipython-input-41-b7aad7986d27>:13: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)
  new_df['tags'] = new_df['tags'].apply(stem)
```

In [42]:
```python
new_df["tags"]
```

Out[42]:
```
0       in the 22nd century, a parapleg marin is dispa...
1       captain barbossa, long believ to be dead, ha c...
2       a cryptic messag from bond' past send him on a...
3       follow the death of district attorney harvey d...
4       john carter is a war-weary, former militari ca...
                              ...
4804    el mariachi just want to play hi guitar and ca...
4805    a newlyw couple' honeymoon is upend by the arr...
4806    "signed, sealed, delivered" introduc a dedic q...
4807    when ambiti new york attorney sam is sent to s...
4808    ever sinc the second grade when he first saw h...
Name: tags, Length: 4806, dtype: object
```

## vectorization

In [50]:
```python
from sklearn.feature_extraction.text import CountVectorizer
cv = CountVectorizer(max_features=5000,stop_words='english')
```

In [48]:
```python
vectors = cv.fit_transform(new_df['tags']).toarray()
```

In [45]:
```python
vectors = cv.fit_transform(new_df['tags']).toarray().shape
```

In [46]:
```python
vectors
```

Out[46]:
```
(4806, 5000)
```

In [51]: `cv.get_feature_names()`

```
---------------------------------------------------------------------------
AttributeError                            Traceback (most recent call last)
<ipython-input-51-631f876094a1> in <module>
----> 1 cv.get_feature_names()

AttributeError: 'CountVectorizer' object has no attribute 'get_feature_names'
```

## model

In [52]: `from sklearn.metrics.pairwise import cosine_similarity`

In [53]: `similarity = cosine_similarity(vectors)`

In [54]: `similarity`

```
Out[54]: array([[1.        , 0.08346223, 0.0860309 , ..., 0.04499213, 0.        ,
         0.        ],
        [0.08346223, 1.        , 0.06063391, ..., 0.02378257, 0.        ,
         0.02615329],
        [0.0860309 , 0.06063391, 1.        , ..., 0.02451452, 0.        ,
         0.        ],
        ...,
        [0.04499213, 0.02378257, 0.02451452, ..., 1.        , 0.03962144,
         0.04229549],
        [0.        , 0.        , 0.        , ..., 0.03962144, 1.        ,
         0.08714204],
        [0.        , 0.02615329, 0.        , ..., 0.04229549, 0.08714204,
         1.        ]])
```

In [55]: `sorted(similarity[0])`

```
Out[55]: [0.0,
 0.0,
 0.0,
 0.0,
 0.0,
 0.0,
 0.0,
 0.0,
 0.0,
 0.0,
 0.0,
 0.0,
 0.0,
 0.0,
 0.0,
 0.0,
 0.0,
 0.0,
```

In [56]: `list(enumerate(similarity[0]))`

```
Out[56]: [(0, 1.0000000000000002),
 (1, 0.08346223261119858),
 (2, 0.08603090020146065),
 (3, 0.0734718358370645),
 (4, 0.18929940971212042),
 (5, 0.10838874619051501),
 (6, 0.04024218182927669),
 (7, 0.14673479641335554),
 (8, 0.0592348877590923),
 (9, 0.0967301666813349),
 (10, 0.1025978352085154l),
 (11, 0.09464970485606021),
 (12, 0.09037128496931669),
 (13, 0.04499212706658476),
 (14, 0.12824729401064427),
 (15, 0.06282808624375433),
 (16, 0.07894736842105264),
 (17, 0.13977653617040256),
 (18, 0.09493290614465533),
```

In [57]:
```python
sorted(list(enumerate(similarity[0])),reverse=True)
## sorting occur's but on the basis of index
```

Out[57]:
```
[(4805, 0.0),
 (4804, 0.0),
 (4803, 0.04499212706658476),
 (4802, 0.046829290579084706),
 (4801, 0.019252140716412975),
 (4800, 0.0),
 (4799, 0.052631578947368425),
 (4798, 0.04223886030955117),
 (4797, 0.0),
 (4796, 0.0),
 (4795, 0.0),
 (4794, 0.0),
 (4793, 0.05407380704358751),
 (4792, 0.0),
 (4791, 0.0),
 (4790, 0.0582716546748065),
 (4789, 0.060833032924035954),
 (4788, 0.0),
 (4787, 0.019672236884115842),
```

In [58]:
```python
sorted(list(enumerate(similarity[0])),reverse=True,key=lambda x:x[1])[1:6]
```

Out[58]:
```
[(1216, 0.28676966733820225),
 (2409, 0.26901379342448517),
 (3730, 0.2605130246476754),
 (507, 0.255608593705383),
 (539, 0.2503866978335957)]
```

In [59]:
```python
def recommend(movie):
    movie_index = new_df[new_df["title"] == movie].index[0]
    distances = similarity[movie_index]
    movies_list = sorted(list(enumerate(distances)),reverse=True,key=lambda x:x[1])[1:6]

    for i in movies_list:
        print(i[0])
```

In [60]:
```python
recommend('Avatar')
```

```
1216
2409
3730
507
539
```

In [61]:
```python
new_df.iloc[1216].title
```

Out[61]:
```
'Aliens vs Predator: Requiem'
```

In [62]:
```python
def recommend(movie):
    movie_index = new_df[new_df["title"] == movie].index[0]
    distances = similarity[movie_index]
    movies_list = sorted(list(enumerate(distances)),reverse=True,key=lambda x:x[1])[1:6]

    for i in movies_list:
        print(new_df.iloc[i[0]].title)
```

In [63]:
```python
recommend('Spectre')
```

```
Quantum of Solace
Skyfall
Never Say Never Again
From Russia with Love
Octopussy
```

In [64]:
```python
recommend('John Carter')
```

```
Riddick
Krrish
The Other Side of Heaven
The Legend of Hercules
Get Carter
```

In [66]:
```python
recommend('The Dark Knight Rises')
```

```
The Dark Knight
Batman Returns
Batman
Batman Forever
Batman Begins
```

In [ ]: